

Optimal IC Task Mapping to Maximize QoS on Heterogeneous Multicore Systems

Lei Mo^{1b}, Member, IEEE, Xinmei Li^{1b}, Student Member, IEEE, Angeliki Kritikakou^{1b}, Member, IEEE, and Pengcheng You^{1b}, Member, IEEE

Abstract—Heterogeneous multicore architectures have become one of the most widely used hardware platforms for embedded systems, where time, energy, and system QoS are the major concerns. The Imprecise Computation (IC) model splits a task into mandatory and optional parts, allowing the trade-off of the above issues. However, existing approaches, to maximize system QoS (Quality-of-Service) under time or energy constraints, use a linear function to model system QoS. Therefore, they become unsuitable for general applications, whose QoS is modeled by a concave function. To deal with this limitation, this brief addresses the Mixed-Integer Non-Linear Programming (MINLP) problem of mapping IC tasks to a set of heterogeneous cores by concurrently deciding which processor executes each task and the number of cycles of optional tasks (i.e., task allocation and scheduling), under real-time and energy supply constraints. Furthermore, as existing solution algorithms either demand high time complexity or only achieve feasible solutions, we propose a novel approach based on problem transformation and dual decomposition that finds an optimal solution while avoiding high computational complexity. Simulation results show that the proposed approach achieves 98% performance of the optimization solver Gurobi, but only with 19.8% of its computation time.

Index Terms—Task mapping, QoS, heterogeneous multicore, imprecise computation, dual decomposition.

I. INTRODUCTION

WITH the increasing requirements for high performance and low task execution delay, multicore platforms have been widely used in various domains, such as Cyber-Physical Systems (CPS). The use of different types of processors in the same platform enables the specialization of the platform

Manuscript received 28 September 2023; accepted 26 October 2023. Date of publication 1 November 2023; date of current version 27 March 2024. This work was supported in part by the National Key Research and Development Program of China under Grant 2022YFF0902800; in part by the Southeast University “Zhishan Scholars” Project under Grant 2242021R40003; in part by the National Natural Science Foundation of China under Grant 61973163, Grant 72201007, Grant 72131001, and Grant 72121002; and in part by the Defense Industrial Technology Development Program under Grant JCKY2020206B068. This brief was recommended by Associate Editor H. Yu. (Corresponding author: Lei Mo.)

Lei Mo and Xinmei Li are with the Key Laboratory of Measurement and Control of CSE, Ministry of Education, School of Automation, Southeast University, Nanjing 210096, China (e-mail: lmo@seu.edu.cn; xinmeili@seu.edu.cn).

Angeliki Kritikakou is with the University of Rennes, INRIA, IRISA, CNRS, 35042 Rennes, France, and also with the Institut Universitaire de France, 75005 Paris, France (e-mail: angeliki.kritikakou@irisa.fr).

Pengcheng You is with the Department of Industrial Engineering and Management, College of Engineering, Peking University, Beijing 100871, China (e-mail: pcyou@pku.edu.cn).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TCSII.2023.3329050>.

Digital Object Identifier 10.1109/TCSII.2023.3329050

into the needs of the application domains, achieving parallel and efficient execution [1]. *Energy efficiency* and *real-time execution* are typically required during the embedded system design because 1) embedded systems have energy constraints, especially when they are powered by the battery, which has limited energy capacity, and 2) real-time responsiveness is required by many critical applications (e.g., target tracking or fire detection) since missing the application’s deadline would lead to serious or even disastrous consequences. In several application domains, such as image processing, robot control, and information gathering, a task can be logically decomposed into a mandatory subtask and an optional subtask [2]. The mandatory subtask must be completed before the deadline, providing a baseline QoS. When it is executed, the optional subtask can increase the QoS, but it is not necessary to be executed completely. Such applications can be modeled by the IC tasks [3]. The way of mapping the IC tasks onto the hardware platforms plays an important role in the final system QoS. This is because the constraints on *energy consumption*, *real-time execution* and *system QoS* often conflict with each other. The longer the optional subtasks are executed, the higher QoS is achieved, while more energy and time are required.

Existing energy-aware task mappings approaches, e.g., [4], [5] aim to minimize energy consumption and typically use a precise task model. Therefore, they do not explore any QoS improvement by adjusting the optional cycles. On the contrary, the QoS-aware task mapping approaches [2], [3], [6], [7], [8], [9], [10] use the IC task model to maximize system QoS under real-time and/or energy constraints. For modeling the system QoS, the most realistic approaches are based on the linear function [2], [3], [6], [7] and the concave function [8], [9], [10]. The concave function is more general since it can characterize more extensive applications than the linear function [11]. However, existing QoS-aware approaches that model QoS through concave functions focus on either uni-processor platforms [10] or with fixed task-to-processor allocation [8] and omit energy-related issues [9]. To address this limitation, we propose a QoS-aware mapping problem for IC tasks with the general concave function on heterogeneous multicore platforms under energy consumption and real-time constraints. Then, we propose a novel polynomial time optimal approach to solve this complex MINLP problem efficiently.

II. SYSTEM MODEL AND PROBLEM FORMULATION

1) *System Model*: We consider a set of periodic and independent IC tasks $\{\tau_1, \dots, \tau_N\}$. They are released at time 0 and share a common hyper-period H [3], [12]. Each IC task τ_i contains a mandatory subtask with M_i cycles and an optional subtask with o_i cycles, where M_i and o_i are measured in

Worst-Case Execution Cycles (WCEC). The optional subtask of τ_i is bounded by $0 \leq o_i \leq O_i$, where O_i is the number of maximum optional cycles. The values of M_i and O_i are determined by the minimum and maximum QoS [13]. We consider a multicore platform with M heterogeneous processors $\{\theta_1, \dots, \theta_M\}$. Each processor θ_k can operate at a voltage/frequency pair (v_k, f_k) . We introduce a parameter $\lambda_{ik} \in (0, 1]$ [5] to describe the execution efficiency of processor θ_k , when θ_k executes task τ_i . Hence, the Worst-Case Execution Time (WCET) of task τ_i , when it is executed at frequency f_k on processor θ_k , is given by $\frac{M_i + o_i}{\lambda_{ik} f_k}$. The processor energy model is adopted from [12], where the power consumption of a processor θ_k is expressed as $P_k^c = P_k^s + P_k^d$. P_k^s is the static power of the processor ready to execute, and P_k^d is the dynamic power of task execution.

2) *Problem Formulation*: Given a set of N IC tasks, the goal is to map these tasks on M heterogeneous processors to maximize the overall system QoS under real-time and energy supply constraints. We consider a concave QoS function $g(\mathbf{o})$. More precisely, we determine 1) which processor the subtask should be executed on (*task allocation*), and 2) how many optional cycles a task should execute (*task scheduling*). To formulate the task mapping problem, we introduce the following variables: 1) define $s \triangleq [s_{ik}]_{N \times M}$ and let binary variable $s_{ik} = 1$, if task τ_i is assigned to processor θ_k , otherwise, $s_{ik} = 0$; and 2) define $\mathbf{o} \triangleq [o_i]_{N \times 1}$ and let continuous variable o_i be the optional cycles of task τ_i . Taking time and energy constraints into account, the Primal Problem (PP) is given by

$$\begin{aligned} \text{PP} : \quad & \min_{s, \mathbf{o}} -g(\mathbf{o}) \\ \text{s.t.} \quad & \mathbf{C}_1 - \mathbf{C}_3, \quad s_{ik} \in \{0, 1\}, \quad 0 \leq o_i \leq O_i. \end{aligned} \quad (1)$$

Since each task τ_i is executed on only one processor, we obtain the following task allocation constraint \mathbf{C}_1 : $\sum_{k \in \mathcal{M}} s_{ik} = 1, \forall i \in \mathcal{N}$, where $\mathcal{N} \triangleq \{1, \dots, N\}$. Considering task allocation decision s_{ik} and task execution cycles $M_i + o_i$, the time spent by processor θ_k to execute all the assigned tasks is known. Thus, the real-time constraint, which restricts all tasks assigned to processor θ_k being executed before the common deadline D , is given by \mathbf{C}_2 : $\sum_{i \in \mathcal{N}} (s_{ik} \frac{M_i + o_i}{\lambda_{ik} f_k}) \leq D, \forall k \in \mathcal{M}$, where $\mathcal{M} \triangleq \{1, \dots, M\}$. Note that the total energy consumed by M processors during the hyper-period H cannot exceed the energy supply E_s . The energy consumption is bounded by \mathbf{C}_3 : $\sum_{k \in \mathcal{M}} (T_k P_k^d + HP_k^s) = \sum_{k \in \mathcal{M}} [\sum_{i \in \mathcal{N}} (s_{ik} \frac{M_i + o_i}{\lambda_{ik} f_k}) P_k^d + HP_k^s] \leq E_s$.

Remark 1: We consider optional cycles o_i as a continuous variable for tractability reasons. When PP is solved, we round the result down. This impact of one cycle is negligible since tasks usually execute hundreds of thousands of cycles [14].

Remark 2: Since binary variable s_{ik} and continuous variable o_i are coupled with each other nonlinearly in \mathbf{C}_2 and \mathbf{C}_3 , and the objective function is concave, PP is a MINLP.

Theorem 1: Task mapping problem (1) is \mathcal{NP} -hard [15].

Due to the page limit, the proofs are omitted here. Although nonlinear item $s_{ik} o_i$ can be linearized by introducing auxiliary variables and adding additional constraints into PP [4], this approach also increases problem size. To circumvent this difficulty, we replace nonlinear item $s_{ik} o_i$ with a new variable t_{ik} , representing the execution time of optional subtask o_i on processor θ_k . Using s_{ik} and t_{ik} , PP can be reformulated as

$$\begin{aligned} \text{PP1} : \quad & \min_{s, t} -g(\mathbf{t}) \\ \text{s.t.} \quad & \mathbf{C}_1, \quad \mathbf{C}'_2, \quad \mathbf{C}'_3, \quad \mathbf{C}_4, \quad s_{ik} \in \{0, 1\}, \quad t_{ik} \geq 0, \end{aligned} \quad (2)$$

where \mathbf{C}'_2 : $\sum_{i \in \mathcal{N}} (s_{ik} \frac{M_i}{\lambda_{ik} f_k} + t_{ik}) \leq D, \forall k \in \mathcal{M}$ and \mathbf{C}'_3 : $\sum_{k \in \mathcal{M}} [\sum_{i \in \mathcal{N}} (s_{ik} \frac{M_i}{\lambda_{ik} f_k} + t_{ik}) P_k^d + HP_k^s] \leq E_s$. With \mathbf{C}_4 : $t_{ik} \leq s_{ik} \frac{O_i}{\lambda_{ik} f_k}, \forall i \in \mathcal{N}, \forall k \in \mathcal{M}$, if a task τ_i is assigned to the processor θ_k , $s_{ik} = 1$. The task execution time of θ_k with respect to the optional subtask of τ_i is bounded by $0 \leq t_{ik} \leq \frac{O_i}{\lambda_{ik} f_k}$, else, $t_{ik} = 0$. Since s_{ik} and t_{ik} are coupled with each other linearly in \mathbf{C}_1 – \mathbf{C}_4 , (2) is an MILP problem.

Lemma 1: Mapping problems (1) and (2) are equivalent.

III. OPTIMAL QoS-AWARE TASK MAPPING METHOD

This section proposes a low computational complexity algorithm to optimally solve the PP1 according to the problem structure. Due to binary variables s , PP1 is a non-convex problem, making this problem difficult to solve. To reduce the computational complexity, we can relax s from binary to continuous variables. However, the feasible region of variables s in PP1 is also changed. Nevertheless, we can achieve an equivalent transformation if we solve the relaxed problem properly, following the three steps below.

1) *Problem Relaxation*: We relax the binary variable s_{ik} to a continuous variable within the range $[0, 1]$. Therefore, (2) is transformed into the following problem:

$$\begin{aligned} \text{PP2} : \quad & \min_{s, t} -g(\mathbf{t}) \\ \text{s.t.} \quad & \mathbf{C}_1 - \mathbf{C}_4, \quad 0 \leq s_{ik} \leq 1, \quad t_{ik} \geq 0. \end{aligned} \quad (3)$$

Let Φ_1^* and Φ_2^* be the optimal objective function values of PP1 and PP2, respectively. Compared with PP1, the feasible region is enlarged in PP2. Since PP1 and PP2 are minimization problems, Φ_1^* is actually an upper bound of Φ_2^* , i.e., $\Phi_1^* \geq \Phi_2^*$.

Theorem 2: The relaxed problem (3) is convex.

Note that the convex problem (3) can be optimally solved by the polynomial-time method. If the optimal solution to the PP2 also satisfies the constraints in the PP1, we get $\Phi_1^* = \Phi_2^*$. To find this type of solution, the constraint $s_{ik} \in \{0, 1\}$ can be combined into the solving process of PP2. Here, we explain how to solve PP1 based on the solution of PP2.

2) *Dual Problem Construction and Refinement*: Instead of solving PP2 directly, we solve its dual problem, as PP2 is a convex problem. To construct the dual problem, we introduce Lagrange multipliers $\alpha, \beta \geq 0, \gamma \geq 0$ and $\mu \geq 0$ with proper dimensions to its constraints. Hence, the Lagrangian is $\mathcal{L}(s, t, \alpha, \beta, \gamma, \mu) = -g(\mathbf{t}) + \sum_{i \in \mathcal{N}} (\sum_{k \in \mathcal{M}} s_{ik} - 1) \alpha_i + \sum_{i \in \mathcal{N}} \sum_{k \in \mathcal{M}} (t_{ik} - s_{ik} \frac{O_i}{\lambda_{ik} f_k}) \beta_{ik} + \sum_{k \in \mathcal{M}} [\sum_{i \in \mathcal{N}} (s_{ik} \frac{M_i}{\lambda_{ik} f_k} + t_{ik}) - D] \gamma_k + \{\sum_{k \in \mathcal{M}} [\sum_{i \in \mathcal{N}} (s_{ik} \frac{M_i}{\lambda_{ik} f_k} + t_{ik}) P_k^d + HP_k^s] - E_s\} \mu$. The dual function $\mathcal{R}(\alpha, \beta, \gamma, \mu)$ is defined as the minimum value of the Lagrangian $\mathcal{L}(s, t, \alpha, \beta, \gamma, \mu)$ over the variables s and t [16]. Thus, for the given α, β, γ and μ , we have objective function $\mathcal{R}(\alpha, \beta, \gamma, \mu) = \min_{s, t} \mathcal{L}(s, t, \alpha, \beta, \gamma, \mu)$, and the Dual Problem (DP) associated with PP2 is given by

$$\begin{aligned} \text{DP} : \quad & \max_{\alpha, \beta, \gamma, \mu} \mathcal{R}(\alpha, \beta, \gamma, \mu) \\ \text{s.t.} \quad & \beta \geq 0, \quad \gamma \geq 0, \quad \mu \geq 0. \end{aligned} \quad (4)$$

From (3), we observe that the inequalities \mathbf{C}_2 – \mathbf{C}_4 are affine, and thus, the Slater's condition [16] is satisfied.

Lemma 2: A strong duality exists between the primal problem (3) and the dual problem (4).

Due to the strong duality, the primal problem and the dual problem have the same optimal solution [16], i.e., solving problem (4) and problem (3) are equivalent. To solve the problem of maximizing $\mathcal{R}(\alpha, \beta, \gamma, \mu)$, we first solve the problem of minimizing $\mathcal{L}(s, t, \alpha, \beta, \gamma, \mu)$. Since variables s and t are coupled with each linearly in $\mathcal{L}(s, t, \alpha, \beta, \gamma, \mu)$, the Lagrangian can be decomposed as $\mathcal{L}(s, t, \alpha, \beta, \gamma, \mu) = \mathcal{L}_1(s, \alpha, \beta, \gamma, \mu) + \mathcal{L}_2(t, \beta, \gamma, \mu) + \mathcal{L}_3(\gamma, \mu)$, where $\mathcal{L}_1(s, \alpha, \beta, \gamma, \mu) = \sum_{i \in \mathcal{N}} \{ \sum_{k \in \mathcal{M}} [\alpha_i - \frac{O_i \beta_{ik} - M_i (\gamma_k + P_k^d \mu)}{\lambda_{i,k} f_k}] s_{ik} - \alpha_i \}$, $\mathcal{L}_2(t, \beta, \gamma, \mu) = \sum_{i \in \mathcal{N}} \sum_{k \in \mathcal{M}} (\beta_{ik} + \gamma_k + P_k^d \mu) t_{ik} - g(t)$, and $\mathcal{L}_3(\gamma, \mu) = \sum_{k \in \mathcal{M}} (HP_k^s \mu - D \gamma_k) - E_s \mu$. Note that $\mathcal{L}_1(s, \alpha, \beta, \gamma, \mu)$ and $\mathcal{L}_2(t, \beta, \gamma, \mu)$ are the functions regarding variables s and t , while $\mathcal{L}_3(\gamma, \mu)$ is fixed under the given multipliers γ and μ . The dual function in (4) is written as $\mathcal{R}(\alpha, \beta, \gamma, \mu) = \min_s \mathcal{L}_1(s, \alpha, \beta, \gamma, \mu) + \min_t \mathcal{L}_2(t, \beta, \gamma, \mu) + \mathcal{L}_3(\gamma, \mu)$. Hence, $\min_s \mathcal{L}_1(s, \alpha, \beta, \gamma, \mu)$ and $\min_t \mathcal{L}_2(t, \beta, \gamma, \mu)$ can be solved separately.

A) To minimize $\mathcal{L}_1(s, \alpha, \beta, \gamma, \mu)$ by s , let $\mathcal{H}_{ik}(\beta, \gamma, \mu) = \frac{O_i \beta_{ik} - M_i (\gamma_k + P_k^d \mu)}{\lambda_{i,k} f_k}$ and substitute it into $\mathcal{L}_1(s, \alpha, \beta, \gamma, \mu)$:

$$\begin{aligned} \min_s \mathcal{L}_1(s, \alpha, \beta, \gamma, \mu) \\ = \min_s \sum_{i \in \mathcal{N}} \{ \sum_{k \in \mathcal{M}} [\alpha_i - \mathcal{H}_{ik}(\beta, \gamma, \mu)] s_{ik} - \alpha_i \}. \end{aligned} \quad (5)$$

Since $[\alpha_i - \mathcal{H}_{i,k}(\beta, \gamma, \mu)] s_{ik}$ is a linear function regarding the variable s_{ik} , (5) can be decomposed into $N \times M$ independent problems with the forms: $\min_{s_{ik}} (\alpha_i - \mathcal{H}_{i,k}(\beta, \gamma, \mu)) s_{ik}, \forall i \in \mathcal{N}, \forall k \in \mathcal{M}$. Comparing α_i with $\mathcal{H}_{i,k}(\beta, \gamma, \mu)$, the task allocation decision s_{ik} can be determined as follows:

$$s_{ik} = \begin{cases} 1 & \text{if } \alpha_i < \mathcal{H}_{i,k}(\beta, \gamma, \mu) \\ [0, 1] & \text{if } \alpha_i = \mathcal{H}_{i,k}(\beta, \gamma, \mu) \\ 0 & \text{if } \alpha_i > \mathcal{H}_{i,k}(\beta, \gamma, \mu) \end{cases} \quad (6)$$

From (6) we have $\min_{s_{ik}} [\alpha_i - \mathcal{H}_{i,k}(\beta, \gamma, \mu)] s_{ik} = [\alpha_i - \mathcal{H}_{i,k}(\beta, \gamma, \mu)]^-$, where $[x]^- \triangleq \min\{0, x\}$. The above equation shows that (5) can be decomposed into N independent functions with the forms: $\sum_{k \in \mathcal{M}} [\alpha_i - \mathcal{H}_{i,k}(\beta, \gamma, \mu)]^-$, $\forall i \in \mathcal{N}$.

Lemma 3: Since each task τ_i is executed on only one processor, to minimize the function $\sum_{k \in \mathcal{M}} [\alpha_i - \mathcal{H}_{i,k}(\beta, \gamma, \mu)]^-$, the optimal task allocation decision s_{ik}^* is determined by:

$$s_{ik}^* = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_k \{\alpha_i - \mathcal{H}_{i,k}(\beta, \gamma, \mu)\} \\ 0 & \text{else} \end{cases} \quad (7)$$

If $\alpha_i - \mathcal{H}_{i,k}(\beta, \gamma, \mu)$ has multiple same minimum values, we randomly select one from these minimum items, as each task is assigned to only one processor.

B) Under the given multipliers β, γ and μ , in order to find an optimal task execution time t^* to minimize $\mathcal{L}_2(t, \beta, \gamma, \mu)$, we solve the following differential equations: $\frac{\partial \mathcal{L}_2(t, \beta, \gamma, \mu)}{\partial t} = 0 \Rightarrow t$. Based on C_2 in PP1, task execution time t^* is influenced by the task allocation decision s^* . The optimal execution time of optimal task t_{ik}^* can be selected as follows:

$$t_{ik}^* = \begin{cases} t_{ik} & \text{if } s_{ik}^* = 1 \\ 0 & \text{else} \end{cases} \quad (8)$$

Thus, we have $\min_t \mathcal{L}_2(t, \beta, \gamma, \mu) = \mathcal{L}_2(t^*, \beta, \gamma, \mu)$, where t^* and $\mathcal{L}_2(t^*, \beta, \gamma, \mu)$ are actually the functions with respect to multipliers β, γ and μ . Substituting (7) and (8) into the objective function of (4), the refined dual function is

Algorithm 1: QoS-Aware Task Mapping Algorithm

```

1 Input: Parameters in (3);
2 Output:  $s$  and  $t$ ;
3 Initialize: Lagrangian multipliers  $\{\alpha, \beta, \gamma, \mu\}$ ,  $j = 0$ ,  $\delta_j = \delta$ ,  $\epsilon$ ;
4 while  $j < N_{\max}$  do
5   Update  $s^{(j)}$  and  $t^{(j)}$  through (7) and (8);
6   Calculate  $g^{(j)}(t)$  and  $\mathcal{L}^{(j)}(s, t, \alpha, \beta, \gamma, \mu)$ ;
7   if  $|\mathcal{L}(s, t, \alpha, \beta, \gamma, \mu) - (-g(t))| < \epsilon$  then
8      $s^* = s^{(j)}$  and  $t^* = t^{(j)}$ ;
9   end
10   $x^{(j+1)} \leftarrow [x^{(j)} + \delta_j \nabla x^{(j)}]^+$ , where  $x = \alpha, \beta, \gamma, \mu$ ;
11  Set  $j \leftarrow j + 1$ , and update  $\delta_{j+1} \leftarrow \delta_j / \sqrt{j}$ ;
12 end

```

$$\begin{aligned} \mathcal{R}(\alpha, \beta, \gamma, \mu) = \mathcal{L}_1(s^*, \alpha, \beta, \gamma, \mu) \\ + \mathcal{L}_2(t^*, \beta, \gamma, \mu) + \mathcal{L}_3(\gamma, \mu). \end{aligned} \quad (9)$$

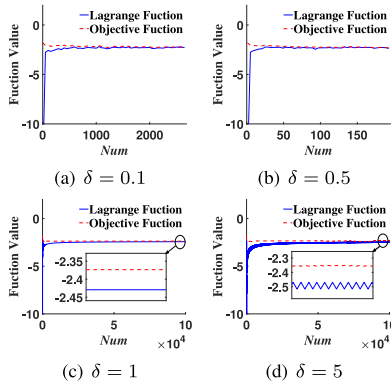
On this basis, we replace the objective function in (4) with the refined dual function (9), and solve the refined dual problem by determining Lagrangian multipliers α, β, γ and μ .

3) *Optimal Task Mapping Solution:* To solve the refined dual problem, an important step is to find the subgradient of the dual function (9). Let $\nabla \alpha_i, \nabla \beta_{ik}, \nabla \gamma_k$ and $\nabla \mu$ denote the subgradient of dual function $\mathcal{R}(\alpha, \beta, \gamma, \mu)$ with respect to the variables $\alpha_i, \beta_{ik}, \gamma_k$ and μ , respectively. Based on the definition of subgradient [17], we have $\nabla \alpha_i = \sum_{k \in \mathcal{M}} s_{i,k}^* - 1$, $\nabla \beta_{i,k} = t_{i,k}^* - s_{i,k}^* \frac{O_i}{\lambda_{i,k} f_k}$, $\nabla \gamma_k = \sum_{i \in \mathcal{N}} (s_{i,k}^* \frac{M_i}{\lambda_{i,k} f_k} + t_{i,k}^*) - D$, and $\nabla \mu = \sum_{k \in \mathcal{M}} [\sum_{i \in \mathcal{N}} (s_{i,k}^* \frac{M_i}{\lambda_{i,k} f_k} + t_{i,k}^*) P_k^d + HP_k^s] - E_s$, where $s_{i,k}^*$ and $t_{i,k}^*$ are corresponding optimal solutions of task allocation and execution time under the given multipliers α, β, γ and μ , respectively. Algorithm 1 summarizes the proposed optimal QoS-aware task mapping method. In each iteration, the Lagrange multipliers are updated with subgradient equations (Line 10), where $[x]^+ \triangleq \max\{0, x\}$. Then, the decisions regarding the task allocation and execution time are updated through (7) and (8) (Line 5). Finally, the iteration stops when $|\mathcal{L}(s, t, \alpha, \beta, \gamma, \mu) - (-g(t))| \leq \epsilon$, where $-g(t)$ is the objective function, $\mathcal{L}(s, t, \alpha, \beta, \gamma, \mu)$ is the Lagrangian, and ϵ is a small positive value that controls the convergence of algorithm [17] (Lines 6–9). The refined dual problem can be solved in polynomial time with the above method. Assume that $\{\alpha^*, \beta^*, \gamma^*, \mu^*\}$ to be the global optimal solution to the refined dual problem. Since the multipliers are optimal, the corresponding solution $\{s^*, t^*\}$ is also optimal.

Remark 3: Transforming dual problem from (2) to (4), we avoid solving s_{ik} directly since s_{ik} can be expressed by other variables α, β, γ and μ . (7) and (8) show that by properly solving the problem of minimizing Lagrangian, we can cut the feasible region of s_{ik} and force s_{ik} to be a binary variable.

Remark 4: (6) and (7) show that s_{ik} can only be set to 0 or 1 when solving the PP2. Moreover, the objective function and the constraints of PP1 and PP2 are the same. Hence, the optimal solution to the PP2 also meets the constraints in the PP1, i.e., $\{s^*, t^*\}$ is also the optimal solution to the PP1.

Remark 5: In each iteration, considering the number of variables, the time complexity is $O(3NM + N + M + 1)$ [18]. In addition, we have $N_{\max} \gg M$ and $N \gg M$, and in the worst case, the maximum number of iterations is N_{\max} . Hence, the time complexity can be simplified as $O(N_{\max} N)$.

Fig. 1. Algorithm convergence under δ varying.

IV. EXPERIMENTAL RESULTS

The simulations are based on a multi-core platform with six heterogeneous processors ($M = 6$). The processor parameters can be found in [4], [12], which are based on 70nm technology. The task efficiency factor is set to $\lambda_{ik} \in (0, 1]$ [5]. The cycles of mandatory subtasks M_i and maximum optional subtasks O_i are assumed within the range $[4 \times 10^7, 6 \times 10^8]$ [12]. This range is calculated from the MiBench and MediaBench benchmarks. To let time and energy constraints change with task number N , we set task deadline $D = H = \theta D_h$ and energy supply $E_s = \eta E_h$, where $\theta, \eta \in (0, 1]$ are time and energy factors, $D_h = \sum_{i \in \mathcal{N}} \max_k \{ \frac{M_i + O_i}{\lambda_{ik} f_k} \}$ and $E_h = \sum_{k \in \mathcal{M}} P_k^s D + \sum_{i \in \mathcal{N}} \max_k \{ \frac{M_i + O_i}{\lambda_{ik} f_k} P_k^d \}$ are the maximum time and energy required to execute N tasks and each task τ_i has $M_i + O_i$ cycles. We adopt a concave QoS function $g(t) = \sum_{i \in \mathcal{N}} \sum_{k \in \mathcal{M}} (-p_{ik} t_{ik} + q_{ik} \sqrt{t_{ik}})$ from [10], where p_{ik} and $q_{ik} \in (0, 1]$ are coefficients related to applications. For Algorithm 1, we set stop criteria $\epsilon = 10^{-3}$, maximum iteration number $N_{\max} = 10^5$, and initial multipliers $\phi^{(0)} = \{\alpha^{(0)}, \beta^{(0)}, \gamma^{(0)}, \mu^{(0)}\} \in [0, 1]$. Although different processor platforms and real application tasks will lead to different PP parameters, the problem structures are still the same, i.e., the proposed method is applicable for different parameter values.

The values of θ and η mainly influence the feasibility of the problem since they change time and energy constraints. In contrast, the values of step-size δ and initial Lagrange multipliers $\phi^{(0)}$ affect the solution quality and the iteration number of our approach when the problem is feasible. In Fig. 1, $N = 20$; we fix the initial Lagrange multipliers $\phi^{(0)}$ and change the step-size δ from 0.1 to 5. The result shows that the algorithm convergence speed increases with δ , as a small step-size will increase the number of iterations and slow the convergence speed ($\delta = 0.1, 0.5$). However, a large step-size will miss the optimal solution ($\delta = 1$) or even lead to oscillation ($\delta = 5$). This result aligns with the stability analysis in [17]. The value of δ should be carefully determined based on the problem parameters. In Fig. 2, we set $N = 20$ and $\delta = 0.5$ and use different initial Lagrange multipliers. *Set 1*: $\phi^{(0)} \in [0, 0.5]$, *Set 2*: $\phi^{(0)} \in [0.5, 1]$, *Set 3*: $\phi^{(0)} \in [1, 1.5]$, and *Set 4*: $\phi^{(0)} \in [1.5, 2]$. When the values of $\phi^{(0)}$ are large (*Set 3* and *Set 4*), the convergence speed is slower, or the result is hard to converge (*Set 4*). However, when the results are stable (*Set 1* and *Set 2*), the influence of $\phi^{(0)}$'s value is limited, compared with δ , as the differences between these sets are usually within 100 iterations. To reduce computation time (iteration number), we can fix $\phi^{(0)}$ and then adjust δ .

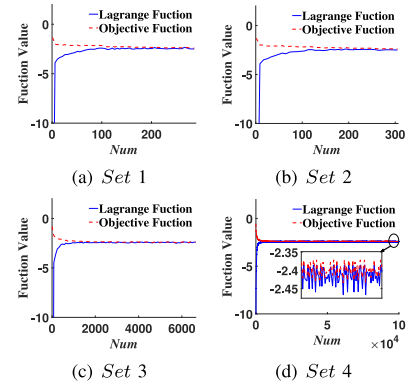
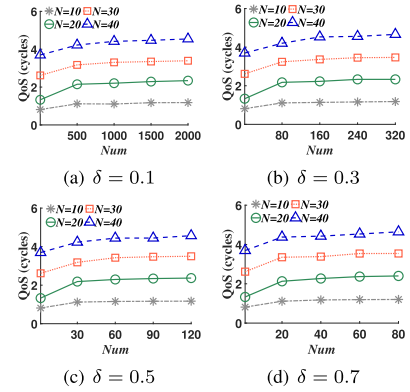
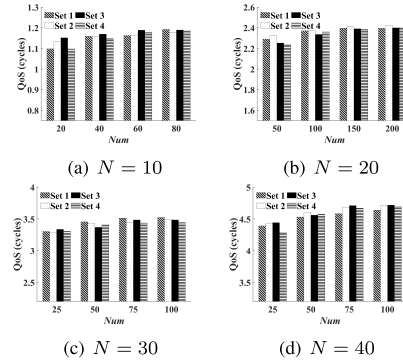
Fig. 2. Algorithm convergence under $\phi^{(0)}$ varying.Fig. 3. QoS of different task sets under δ varying.Fig. 4. QoS of different task sets under $\phi^{(0)}$ varying.

Fig. 3 and Fig. 4 compare system QoS achieved by the proposed approach under different iteration numbers Num and task numbers N . In Fig. 3, we fix $\phi^{(0)}$ and change δ , while in Fig. 4, we fix $\delta = 0.5$ and use different $\phi^{(0)}$. Note that Fig. 1 and Fig. 2 show the results when the iterations are stopped (i.e., the convergence criteria are satisfied); Fig. 3 and Fig. 4 show the results achieved by a given iteration number. Since our approach is iteration-based, the more iterations are performed, the better the solution quality, as shown in Fig. 3. Our approach initially converges very fast. With the iteration number increasing, the improvement in solution quality is reduced. For instance, when $N = 10$ and $\delta = 0.1$, the QoS loss between $Num = 500$ and $Num = 2000$ is less than 5%. However, the computation time of $Num = 2000$ is four times of $Num = 500$. This is because the adaptive step-size method is used in Algorithm 1. The step-size is updated by

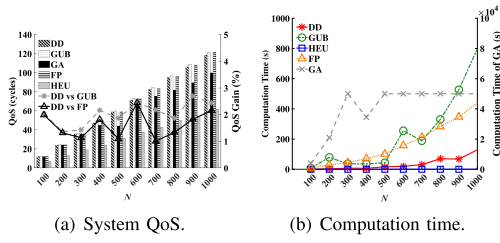


Fig. 5. QoS and computation time comparisons of different methods.

$\delta_{j+1} \leftarrow \delta_j / \sqrt{j}$, i.e., δ_{j+1} decreases with iteration number j . In some time-sensitive applications, the iteration can be stopped when solution quality already reaches the requirements (e.g., the time and energy constraints are met or the gap between the achieved and optimal solutions is within a threshold), even if the convergence criteria are not satisfied. Fig. 4 shows that under the given N and Num parameters, the QoS differences between different initial multipliers $\phi^{(0)}$ are small, as δ plays a more important role than $\phi^{(0)}$ to influence solution quality.

Fig. 5 compares system performance (QoS and computation time) of the proposed approach, i.e., Dual Decomposition-based (DD) approach, with the state-of-the-art optimization solver Gurobi (GUB), Genetic Algorithm (GA), Feasibility Pump method (FP) [19] and a two-step heuristic (HEU) [12]. The MATLAB optimization toolbox provides GA. The QoS gain between GUB and DD is defined as $\frac{f_{GUB}(N) - f_{DD}(N)}{f_{GUB}(N)}$, where $f_{GUB}(N)$ and $f_{DD}(N)$ is the QoS achieved by GUB and DD under the given N parameter. We set $\delta = 0.5$ and vary task number N from 100 to 1000, and each time $\phi^{(0)}$ is randomly selected within $[0, 1]$. Fig. 5(a) shows that DD outperforms GA and HEU in terms of QoS improvement (7.8% and 48% on average). Although GA can solve complex mixed programming problems, such as MINLP and MILP, the solution's optimality is hard to guarantee. The HEU method, which determines task allocation and optional cycle adjustment in sequence, has a lower system QoS since separating task allocation and optional cycle adjustment reduces the feasibility region of the original problem. To enhance solution quality, the subproblems should be considered simultaneously. Fig. 5(b) compares the computation time of different methods. DD achieves about 98% and 98.3% (on average) performance of GUB and FP, but only with 19.8% and 14.4% computation time. This is because DD (1) is solved in a convex manner, and the computation time is less than solving a MINLP. The result shows that the computation time of DD almost grows linearly with task number N . This result aligns with the complexity analysis. Since our approach can transform MINLP into a convex problem equivalently, it achieves a balance between computation time and solution quality. Although DD and GA are both iteration-based methods, the GA structure is more complex than DD, as it generates new populations through several procedures at each iteration, e.g., selection, reproduction, mutation and crossover. Fig. 5 shows that when N is large, DD obtains a better QoS in a short time compared with GA. For the iteration approach, only when the step-size is small enough and the iteration number is large enough, we can reach the optimal solution [17]. The above results of DD can be further improved through better δ and Num parameters.

V. CONCLUSION

This brief studied the QoS-aware task mapping problem for IC tasks on heterogeneous multicore platforms to maximize

the system QoS under real-time and energy constraints, formulated as an MINLP problem. Based on the problem transformation and dual decomposition, we proposed a novel method to efficiently solve this problem in low computation time, whose solution is the same as the optimal solution to the original problem. Finally, the simulation results show the effectiveness of the proposed method, which outperforms other methods in terms of QoS improvement and computation time.

REFERENCES

- [1] J. Zhou, Y. Shen, L. Li, C. Zhuo, and M. Chen, "Swarm intelligence-based task scheduling for enhancing security for IoT devices," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 42, no. 6, pp. 1756–1769, Jun. 2023.
- [2] S. Chakraborty, S. Saha, M. Sjalander, and K. McDonald-Maier, "PREPARE: Power-aware approximate real-time task scheduling for energy-adaptive QoS maximization," *ACM Trans. Embed. Comput. Syst.*, vol. 20, no. 5, pp. 1–25, Sep. 2021.
- [3] L. Mo, A. Kritikakou, and O. Sentieys, "Decomposed task mapping to maximize QoS in energy-constrained real-time multicores," in *Proc. IEEE Int. Conf. Comput. Design*, 2017, pp. 493–500.
- [4] G. Chen, K. Huang, and A. Knoll, "Energy optimization for real-time multiprocessor system-on-chip with optimal DVFS and DPM combination," *ACM Trans. Embed. Comput. Syst.*, vol. 13, no. 3, pp. 1–21, Mar. 2014.
- [5] D. Li and J. Wu, "Minimizing energy consumption for frame-based tasks on heterogeneous multiprocessor platforms," *IEEE Trans. Parallel Distrib. Syst.*, vol. 26, no. 3, pp. 810–823, Mar. 2015.
- [6] M. Ansari et al., "Thermal-aware standby-sparing technique on heterogeneous real-time embedded systems," *IEEE Trans. Emerg. Topics Comput.*, vol. 10, no. 4, pp. 1883–1897, Oct.–Dec. 2022.
- [7] K. Cao, J. Zhou, G. Xu, T. Wei, and S. Hu, "Exploring renewable-adaptive computation offloading for hierarchical QoS optimization in fog computing," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 39, no. 10, pp. 2095–2108, Oct. 2020.
- [8] H. Yu, Y. Ha, B. Veeravalli, F. Chen, and H. El-Sayed, "DVFS-based quality maximization for adaptive applications with diminishing return," *IEEE Trans. Comput.*, vol. 70, no. 5, pp. 803–816, May 2021.
- [9] G. M. Tchamgoue, K. H. Kim, Y. K. Jun, and W. Y. Lee, "Compositional real-time scheduling framework for periodic reward-based task model," *J. Syst. Softw.*, vol. 86, no. 6, pp. 1712–1724, Jun. 2013.
- [10] L. A. Cortes, P. Eles, and Z. Peng, "Quasi-static assignment of voltages and optional cycles in imprecise-computation systems with energy considerations," *IEEE Trans. Very Large Scale Integr. Syst.*, vol. 14, no. 10, pp. 1117–1129, Oct. 2006.
- [11] H. Aydin, R. Melhem, D. Mosse, and P. Mejia-Alvarez, "Optimal reward-based scheduling for periodic real-time tasks," *IEEE Trans. Comput.*, vol. 50, no. 2, pp. 111–130, Feb. 2001.
- [12] J. Zhou, J. Yan, T. Wei, M. Chen, and X. S. Hu, "Energy-adaptive scheduling of imprecise computation tasks for QoS optimization in real-time MPSoC systems," in *Proc. Design, Autom. Test Eur.*, 2017, pp. 1402–1407.
- [13] H. Yu, Y. Ha, and B. Veeravalli, "Quality-driven dynamic scheduling for real-time adaptive applications on multiprocessor systems," *IEEE Trans. Comput.*, vol. 62, no. 10, pp. 2026–2040, Oct. 2013.
- [14] L. Mo, A. Kritikakou, and O. Sentieys, "Energy-quality-time optimized task mapping on DVFS-enabled multicores," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 37, no. 11, pp. 2428–2439, Nov. 2018.
- [15] S. Burer and A. N. Letchford, "Non-convex mixed-integer nonlinear programming: A survey," *Surv. Oper. Res. Manage. Sci.*, vol. 17, no. 2, pp. 97–106, Jul. 2012.
- [16] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K: Cambridge Univ. Press, 2004.
- [17] S. Boyd, L. Xiao, and A. Mutapcic, "Subgradient methods," Lecture Notes EE392o, Stanford Univ., Stanford, CA, USA, 2003.
- [18] O. He, S. Dong, W. Jang, J. Bian, and D. Z. Pan, "UNISM: Unified scheduling and mapping for general networks on chip," *IEEE Trans. Very Large Scale Integr. Syst.*, vol. 20, no. 8, pp. 1496–1509, Aug. 2012.
- [19] D. E. Bernal, S. Vigerske, F. Trespalcacios, and I. E. Grossmann, "Improving the performance of DICOPT in convex MINLP problems using a feasibility pump," *Optim. Methods Softw.*, vol. 35, no. 1, pp. 171–190, 2020.